# Deloitte.

# Webinar ''ChatGPT in business''

🕐 **May 15, 2023, 17:00**

💬 **English**

📍 **Zoom**

## Agenda:

**17:05 – 17:20 | The slow path to intelligence – the current state and future prospects of large language models**
**Dr. Ēvalds Urtāns,** asya.ai

**17:20 – 17:35 | Where and when to implement Chat GPT for business problems**
**Igor Rodin**, Deloitte CE

**17:35 – 17:50 | Practical guidance on building GPT-4 based applications in business**
**Dr. Romāns Taranovs**, Deloitte CE

# Deloitte webinar
# The slow path to intelligence

**Evalds Urtans**
asya.ai
CEO

# What is AI?

# What is AI?



$$i_t = \sigma(W_i * [\mathcal{X}_t, \mathcal{H}_{t-1}] + b_i)$$

$$f_t = \sigma(W_f * [\mathcal{X}_t, \mathcal{H}_{t-1}] + b_f)$$

$$e_{t,z} = V_e \cdot \tanh(W_e * [\mathcal{X}_{t,z}, \mathcal{H}_{t-1}] + b_e)$$

$$\alpha_{t,z} = \frac{\exp(e_{t,z})}{\sum_{j=1}^{\tau} \exp(e_{t,j})}$$

$$p_t = \sum_{j=1}^{\tau} \alpha_{t,j} \tilde{\mathcal{X}}_{t,j}$$

$$n_t = \sigma(W_n * [\mathcal{X}_t, \mathcal{H}_{t-1}] + b_n)$$

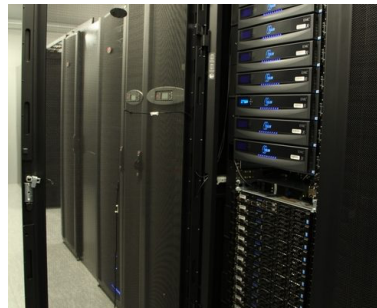$$g_t = \tanh(W_g * [p_t, \mathcal{H}_{t-1}] + b_g)$$

$$\mathcal{C}_t = f_t \circ \mathcal{C}_{t-1} + i_t \circ a_t + n_t \circ g_t$$

$$a_t = \tanh(W_a * [\mathcal{X}_t, \mathcal{H}_{t-1}] + b_a)$$

$$o_t = \sigma(W_o * [\mathcal{X}_t, \mathcal{H}_{t-1}] + b_o)$$

$$\mathcal{H}_t = o_t \circ \tanh(\mathcal{C}_t)$$

# What is AI?

- **Linear algebra**
- **Calculus**
- **Probability theory**
- **Information theory**
- **10% programming**

$$i_t = \sigma(W_i * [\mathcal{X}_t, \mathcal{H}_{t-1}] + b_i)$$

$$f_t = \sigma(W_f * [\mathcal{X}_t, \mathcal{H}_{t-1}] + b_f)$$

$$e_{t,z} = V_e \cdot \tanh(W_e * [\mathcal{X}_{t,z}, \mathcal{H}_{t-1}] + b_e)$$

$$\alpha_{t,z} = \frac{\exp(e_{t,z})}{\sum_{j=1}^{\tau} \exp(e_{t,j})}$$

$$p_t = \sum_{j=1}^{\tau} \alpha_{t,j} \tilde{\mathcal{X}}_{t,j}$$

$$n_t = \sigma(W_n * [\mathcal{X}_t, \mathcal{H}_{t-1}] + b_n)$$

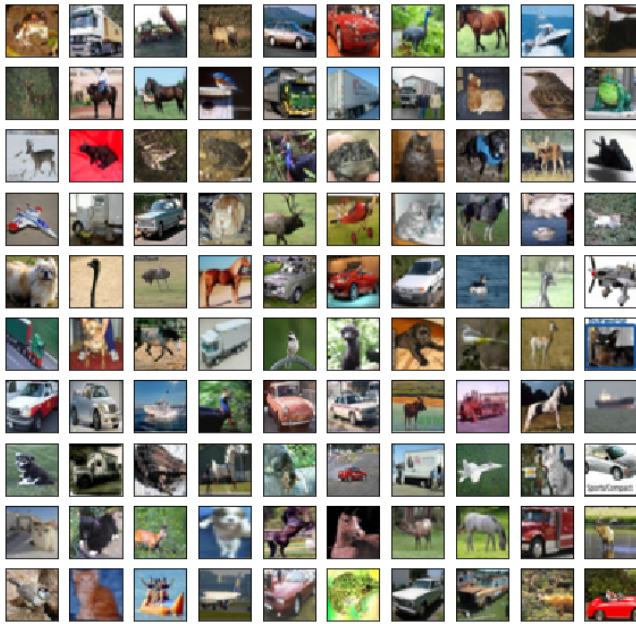$$g_t = \tanh(W_g * [p_t, \mathcal{H}_{t-1}] + b_g)$$

$$\mathcal{C}_t = f_t \circ \mathcal{C}_{t-1} + i_t \circ a_t + n_t \circ g_t$$

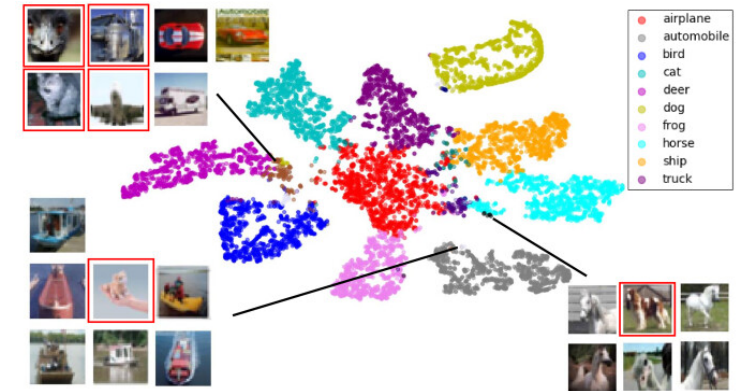$$a_t = \tanh(W_a * [\mathcal{X}_t, \mathcal{H}_{t-1}] + b_a)$$

$$o_t = \sigma(W_o * [\mathcal{X}_t, \mathcal{H}_{t-1}] + b_o)$$

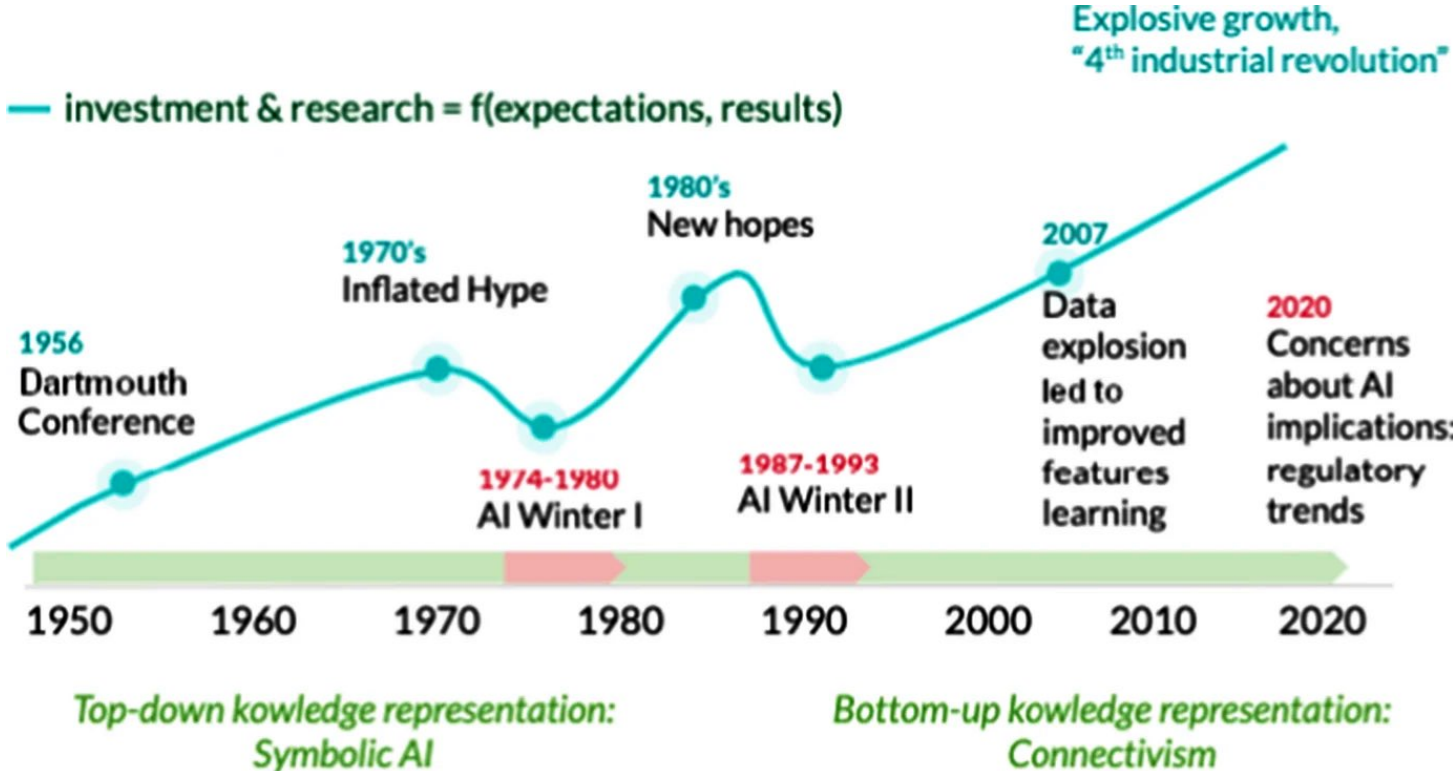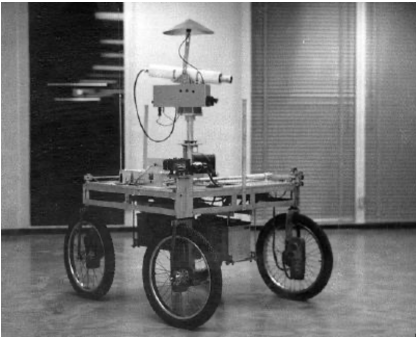$$\mathcal{H}_t = o_t \circ \tanh(\mathcal{C}_t)$$

# What is AI?



$$i_t = \sigma(W_i * [\mathcal{X}_t, \mathcal{H}_{t-1}] + b_i)$$

$$f_t = \sigma(W_f * [\mathcal{X}_t, \mathcal{H}_{t-1}] + b_f)$$

$$e_{t,z} = V_e \cdot \tanh(W_e * [\mathcal{X}_{t,z}, \mathcal{H}_{t-1}] + b_e)$$

$$\alpha_{t,z} = \frac{\exp(e_{t,z})}{\sum_{j=1}^{\tau} \exp(e_{t,j})}$$

$$p_t = \sum_{j=1}^{\tau} \alpha_{t,j} \tilde{\mathcal{X}}_{t,j}$$

$$n_t = \sigma(W_n * [\mathcal{X}_t, \mathcal{H}_{t-1}] + b_n)$$

$$g_t = \tanh(W_g * [p_t, \mathcal{H}_{t-1}] + b_g)$$

$$\mathcal{C}_t = f_t \circ \mathcal{C}_{t-1} + i_t \circ a_t + n_t \circ g_t$$

$$a_t = \tanh(W_a * [\mathcal{X}_t, \mathcal{H}_{t-1}] + b_a)$$

$$o_t = \sigma(W_o * [\mathcal{X}_t, \mathcal{H}_{t-1}] + b_o)$$

$$\mathcal{H}_t = o_t \circ \tanh(\mathcal{C}_t)$$

# History



investment & research = f(expectations, results)

Explosive growth, "4th industrial revolution"

1980's New hopes

1970's Inflated Hype

2007 Data explosion led to improved features learning

2020 Concerns about AI implications: regulatory trends

1956 Dartmouth Conference

1974-1980 AI Winter I

1987-1993 AI Winter II

1950    1960    1970    1980    1990    2000    2010    2020

Top-down kowledge representation: Symbolic AI

Bottom-up kowledge representation: Connectivism

**Hans Moravec's Robots, 1975**



**Tesla FSD, 2023**

# Time to Reach 100M Users

## Months to get to 100 million global Monthly Active Users

- Google Translate: 78
- Uber: 70
- Telegram: 61
- Spotify: 55
- Pinterest: 41
- Instagram: 30
- TikTok: 9
- ChatGPT: 2

@EconomyApp

APP ECONOMY INSIGHTS

# Language modelling, Tokens

The color of the sky is → Language Model → blue

91%

Don't waste food

▼

**Subword Tokenization**

▼

| Do | n't | waste | food |

# Old way - RNN, LSTM



- S. Hochreiter, J. Schmidhuber 1995
- Not-parellizable
- Limited memory capacity
- Small VRAM footprint
- Weaker performance: 74% acc. vs 82% acc (New way)

# New way - Transformer



Transformers at school  Transformers at college  Transformers today

# New way - Transformer

The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
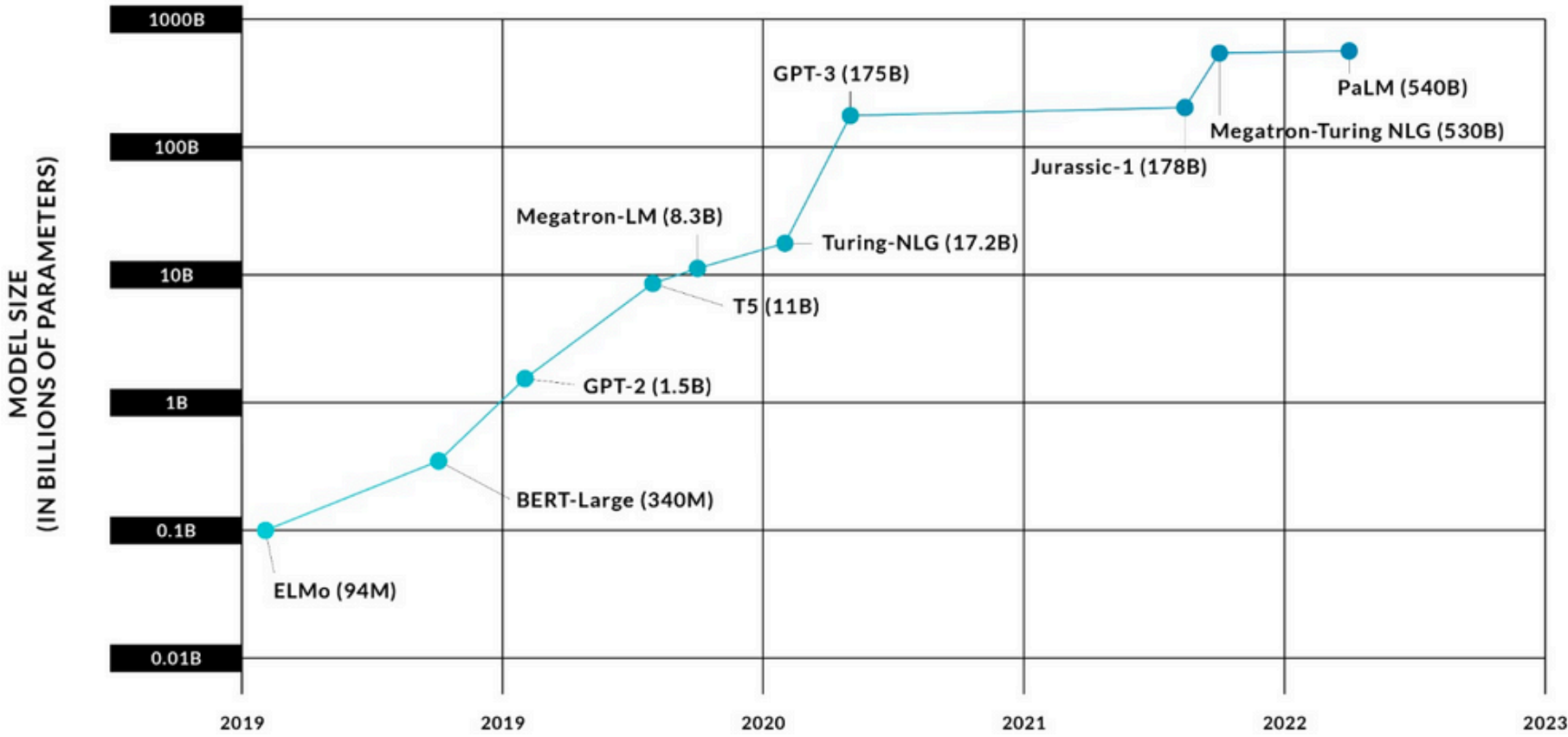The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .

- BERT (Google), GPT (OpenAI), 2018
- Transformer architecture
- Parellizable
- No memory*
- Very large VRAM footprint
- Limited context length ~2048 tokens** (or tradeoffs)
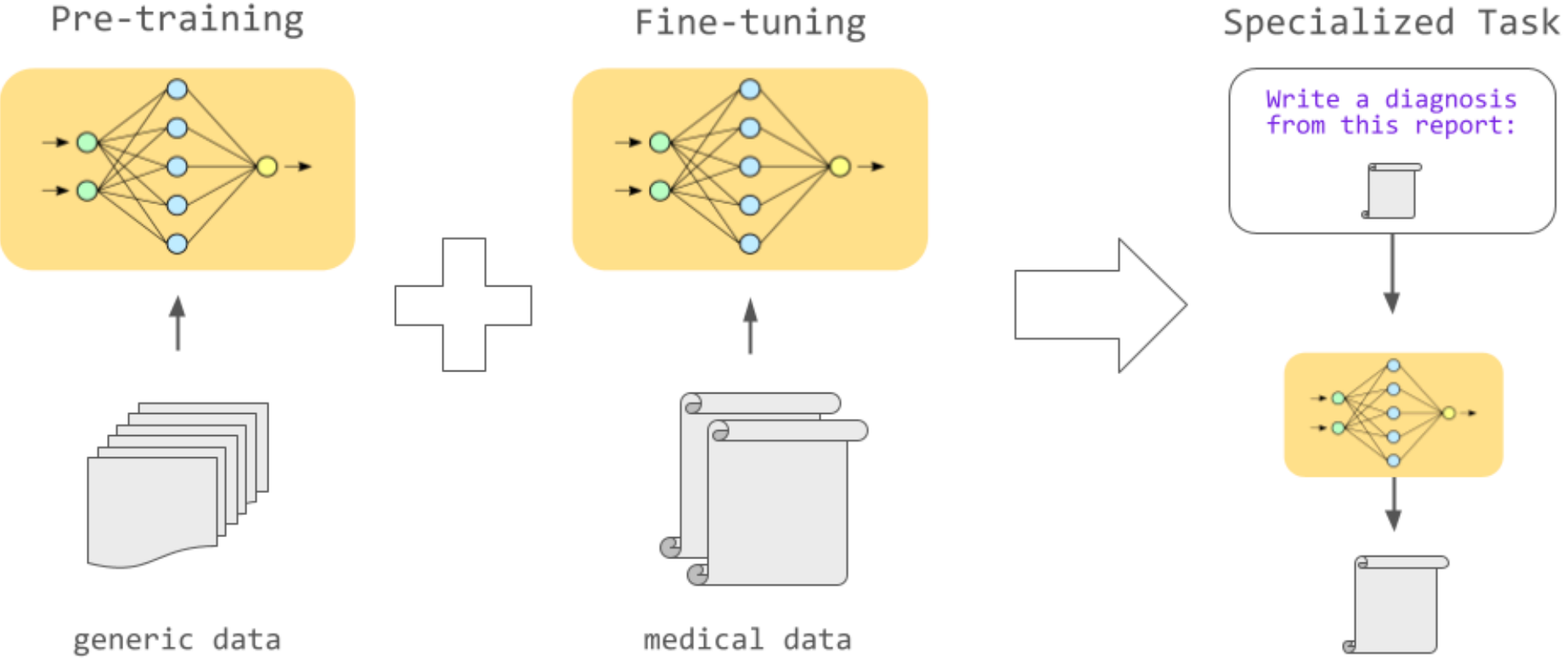- CommonCrawl (10 years), 410b tokens
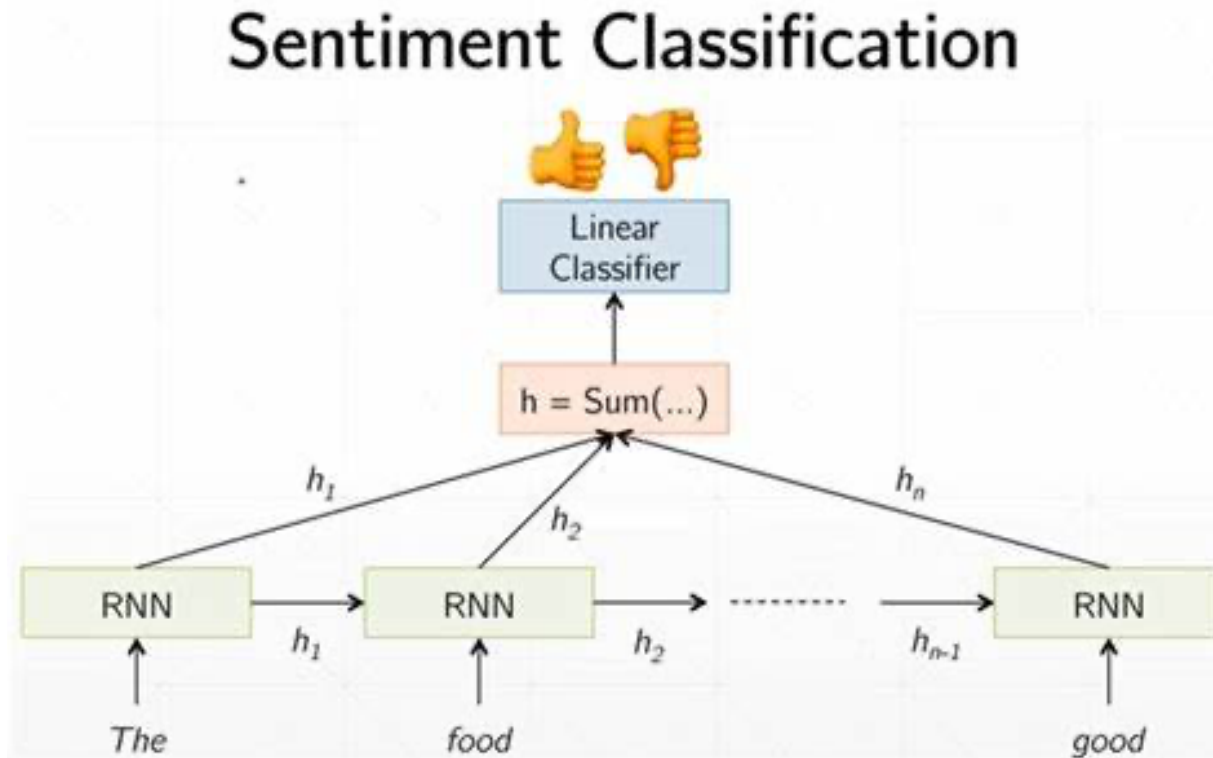
# New way - Transformer

# Use-cases - Train from scratch



Language Model Sizes Over Time

# Use-cases - Fine tunning

# Use-cases - Auxilary tasks



Sentiment Classification

- Classification, Regression, segmentation

- Need labelled data at least 10k text input

- Expect higher accuracy 90%+

# Text embeddings
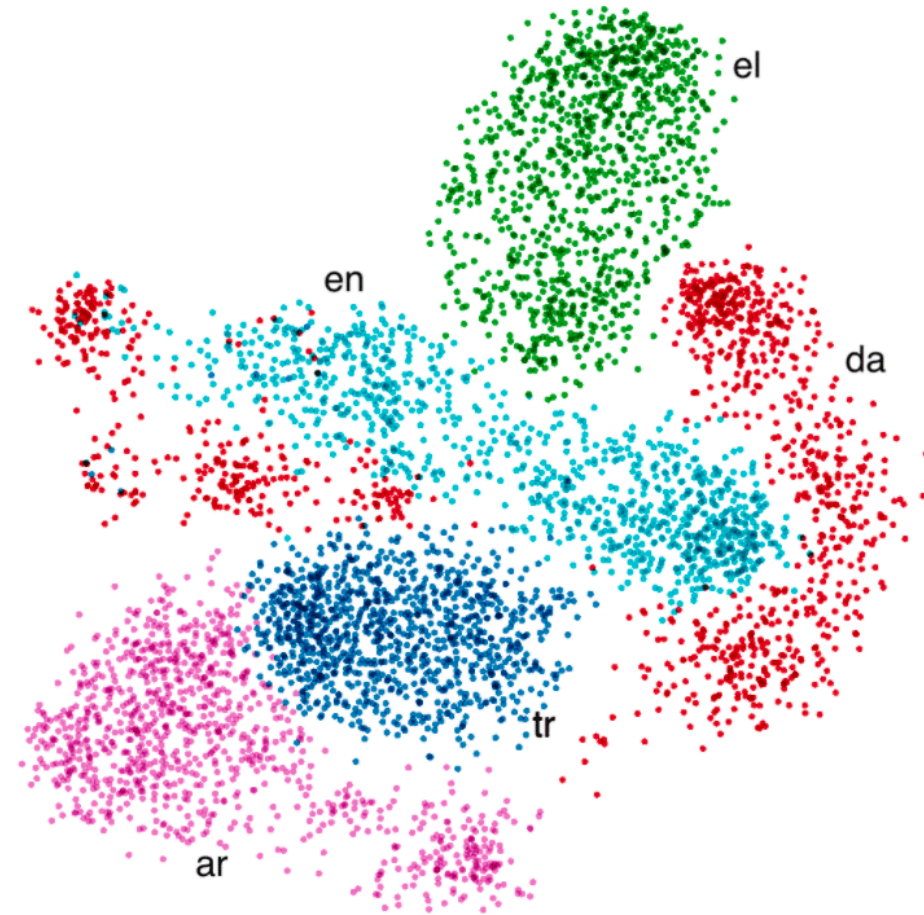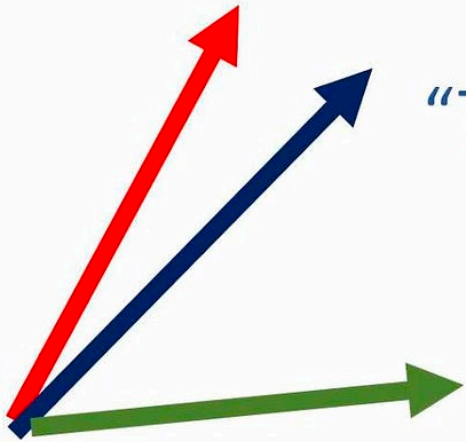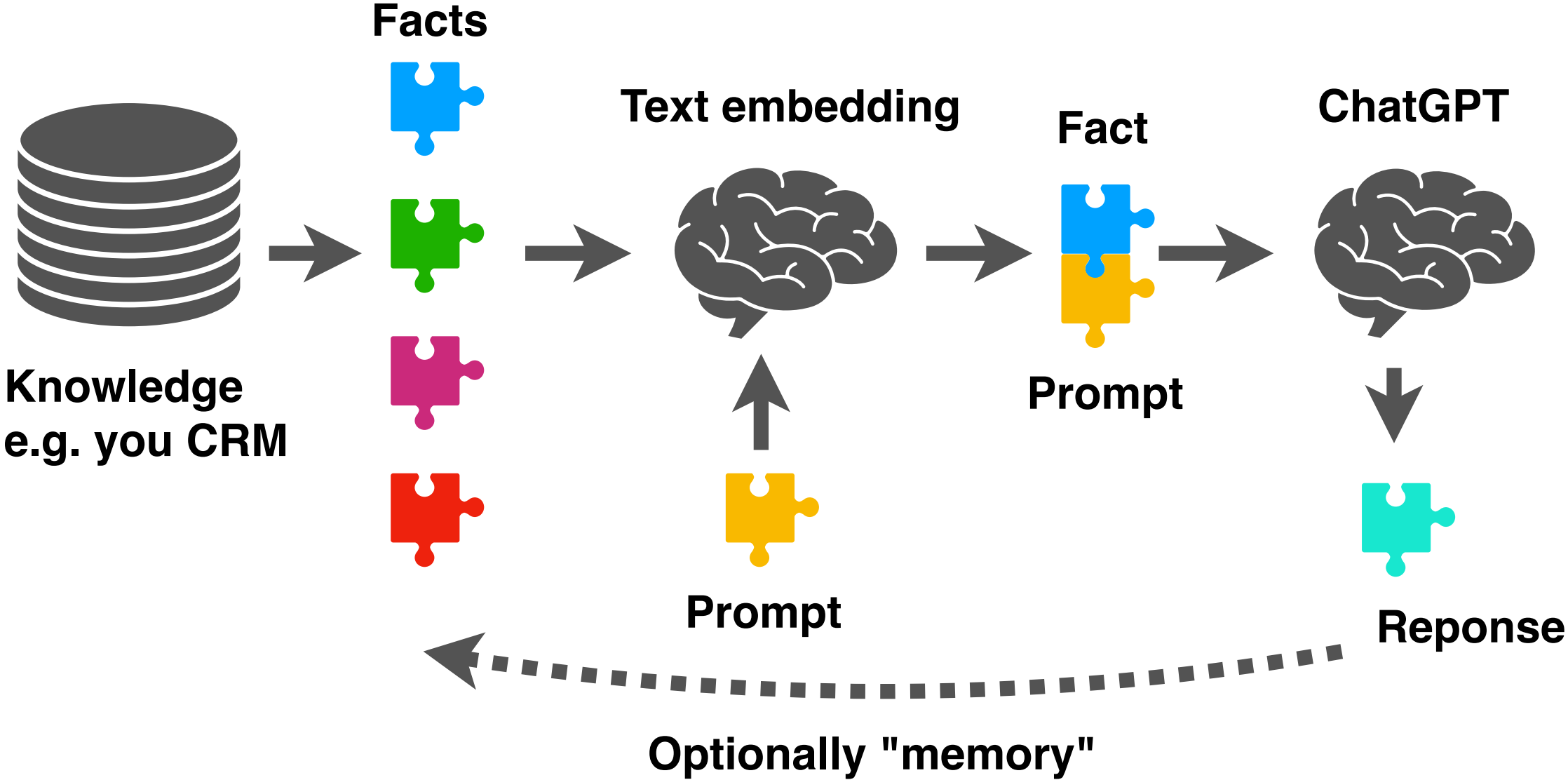
# Use-cases - Zero-shot + Knowledge base

**Facts**

**Text embedding**

**Fact**

**ChatGPT**

**Knowledge e.g. you CRM**

**Prompt**

**Prompt**

**Reponse**

**Optionally "memory"**

# Hardware resources

**Necessary hardware:**
- **AWS EC2, Google Cloud, Oracle cloud - 10k EUR/mon**
- **RTU HPC - <10k EUR/mon**
- **nVidia GPU V100, A100, A10 40GB (min.) - 20k EUR/per unit**
- **Google Collab - 50 EUR/mon**
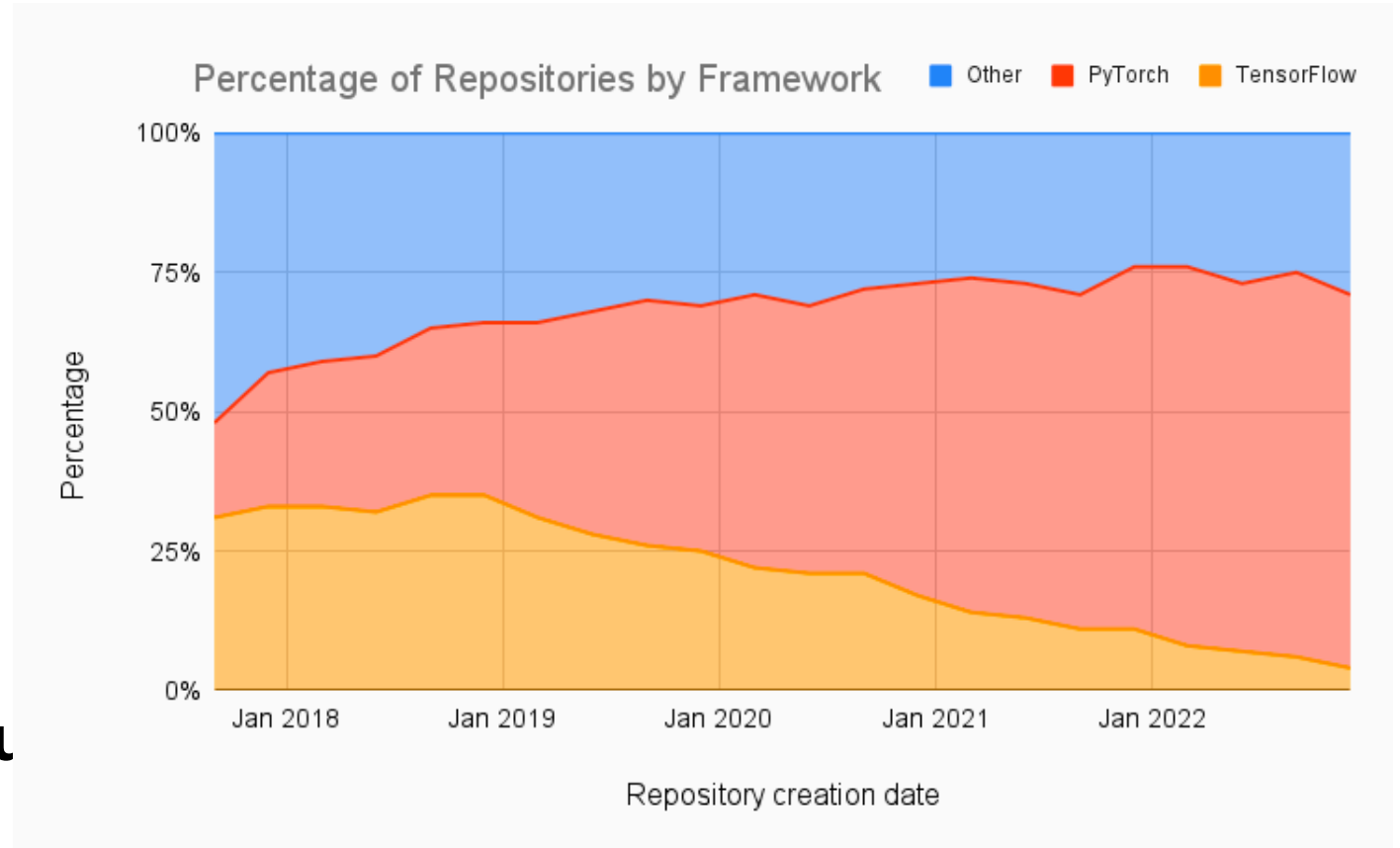
# Software resources

**Best Open-Source LLMs**
- **FLAN-T5 XXL**
- **GPT-JT**
- **Bloom**
- **Open Assistant**
- **LLaMA\*\***

**Model sources:**
- **huggingface.io**
- **torchvision, torchtext, torchhu**

**Model programming:**
- **PyTorch**
- **ONNX (cross-platform deployment)**



Percentage of Repositories by Framework

# Tools to explore

**Search engines:**
- **perplexity.ai**
- **chat.you.com**

**Productivity:**
- **chatpdf.com**

**Content:**
- **jasper.ai**
- **writesonic.ai**



Google traffic share vs Bing + chatGPT

chatGPT Launches

Source: ARK Invest, SimilarWeb

# Deloitte.

# Deloitte webinar
# Chat GPT for business problems

**Evalds Urtans**
evalds@asya.ai
CEO